

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

Applicant: Takahiko KAWATANI
Serial No.: NEW
Filed: March 4, 2004
For: DOCUMENT AND PATTERN CLUSTERING METHOD AND APPARATUS

CLAIM FOR PRIORITY AND SUBMISSION OF PRIORITY DOCUMENT

Commissioner for Patents
P. O. Box 1450
Alexandria, VA 22313-1450

March 4, 2004

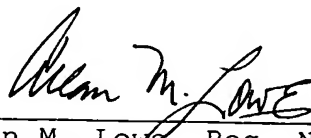
Sir:

In accordance with the provisions of 35 U.S.C. §119, Applicant hereby claims the benefit of the filing dates of Japanese Patent Application Nos. 2003-105867, filed March 5, 2003, and 2004-30629, filed February 6, 2004. A certified copy of JP 2003-105867 is attached; a certified copy of JP 2004-30629 will be filed in due course. The Examiner is respectfully requested to acknowledge Applicant's claim for priority, as well as receipt of the certified copy of the priority document.

Respectfully submitted,

LOWE HAUPTMAN GILMAN & BERNER, LLP

By:


Allan M. Lowe, Reg. No. 19,641

1700 Diagonal Road, Suite 300
Alexandria, VA 22314
703-684-1111 telephone
703-518-5499 telecopier
AML:rk

日 本 国 特 許 庁
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日 2 0 0 3 年 3 月 5 日
Date of Application:

出 願 番 号 特 願 2 0 0 3 - 1 0 5 8 6 7
Application Number:

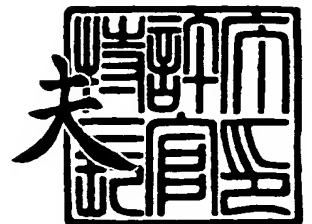
[ST. 10/C]: [J P 2 0 0 3 - 1 0 5 8 6 7]

出 願 人 ヒューレット・パカード・カンパニー
Applicant(s):

2 0 0 3 年 7 月 2 3 日

特許庁長官
Commissioner,
Japan Patent Office

今 井 康 夫



出証番号 出証特 2 0 0 3 - 3 0 5 8 6 2 0

【書類名】 特許願

【提出日】 平成15年 3月 5日

【整理番号】 200309904

【あて先】 特許庁長官 殿

【国際特許分類】 G06F 17/21

【発明者】

【住所又は居所】 東京都杉並区高井戸東3丁目29番21号 日本ヒュー
レット・パッカード株式会社内

【氏名】 川谷 隆彦

【特許出願人】

【識別番号】 398038580

【氏名又は名称】 ヒューレット・パッカード・カンパニー

【代理人】

【識別番号】 100082946

【弁理士】

【氏名又は名称】 大西 昭広

【手数料の表示】

【予納台帳番号】 061492

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 0300115

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 文書のクラスタリング方法及び装置

【特許請求の範囲】

【請求項 1】

以下の (a) から (e) のステップを有する、一つまたは複数の文書セグメントを持つ複数の文書から成る文書集合に対して、文書間の関係に基づき文書をグループ化するクラスタリング方法、

(a) 前記文書集合と共通性の高い文書を前記文書集合から選択しクラスターとするステップ、

(b) 前記選択された文書と類似性の高い全ての文書を前記クラスターに追加するステップ、

(c) 前記クラスターに存在する文書と共通性の高い全ての文書を前記クラスターに繰返し追加するステップ、

(d) 既に求められたクラスターと共通性の高い文書を除外した文書集合を前記文書集合として前記 (a)、(b) 及び (c) のステップを繰返すステップ、

(e) 前記クラスターが重複する場合に、冗長なクラスターを除去するステップ。

【請求項 2】

前記 (a) のステップは、さらに以下の (a-1) から (a-4) のステップを含むことを特徴とする請求項 1 に記載の方法、

(a-1) 前記文書セグメント毎に、前記文書セグメントに出現する用語に対応する成分の値を 1、他の値は 0 とする文書セグメントベクトルを生成するステップと、

(a-2) 文書セグメントベクトルより生成する共起行列を前記文書集合の各文書に対して求めるステップと、

(a-3) 前記文書集合に含まれる全文書から文書頻度行列を求めるステップと、

(a-4) 前記共起行列及び文書頻度行列を用いて、前記文書集合と共通性の高い文書を選択するステップ。

【請求項 3】

前記 (c) のステップは、さらに以下の (c-1) から (c-3) のステップを含むことを特徴とする請求項 1 に記載の方法、

(c-1) 各用語の入力文書集合における出現文書数と前記クラスター内での出現文書数から文書頻度比を求めるステップ、

(c-2) 前記クラスター内での出現文書数や前記文書頻度比の情報をを用いて用語や用語対を選択し、選択された用語や用語対を用いて文書共通度を求めるステップ、

(c-3) 前記文書共通度に基づいて、前記前記選択された文書と共通性の高い文書を選択するステップ。

【請求項 4】

前記 (d) のステップは、さらに以下の (d-1) から (d-4) のステップを含むことを特徴とする請求項 1 に記載の方法、

(d-1) 各文書と各クラスターとの文書共通度を求めるステップ、

(d-2) 各クラスター毎に、前記文書共通度が、着目クラスターに対しては所定値よりも大きく、他の全てのクラスターに対して所定値よりも小さい文書数を求めるステップ、

(d-3) 前記文書数が所定値よりも小さいクラスターが存在する場合には、最少の前期文書数を有するクラスターを削除するステップ、

(d-4) 前記 (d-1) から (d-3) を繰返すステップ。

【請求項 5】

以下の (a) から (1) のステップを有する、一つまたは複数の文書セグメントを持つ複数の文書から成る文書集合に対して、同じ話題を有する文書をグループ化するクラスタリング方法、

(a) 前記文書セグメント毎に、前記文書セグメントに出現する用語に対応する成分の値を 1、他の値は 0 とする文書セグメントベクトルを生成するステップと

(b) 文書セグメントベクトルより生成する共起行列を前記文書集合の各文書に対して求めるステップと、

- (c) 前記文書集合に含まれる全文書から文書頻度行列を求めるステップと、
- (d) その時点で求められているどのクラスターに対しても文書共通度が一定値以下の文書の集合に対して共通共起行列を求めるステップと、
- (e) 前記その時点で求められているどのクラスターに対しても文書共通度が一定値以下の文書の集合の中から (d) で求められた共通共起行列を用いてクラスターの種となる文書を抽出し、クラスターの種文書の近隣文書からなる初期クラスターを作成するステップと、
- (f) その時点で着目クラスターに一時的に帰属する文書集合から共通共起行列と文書頻度行列を求めるステップと、
- (g) (c) 及び (f) で求められる文書頻度行列を比較して、各用語及び用語対の着目クラスターに対する特有さの情報を求めるステップと、
- (h) 各文書の着目クラスターに対する文書共通度を、(f) で求められる共通共起行列と、(g) で求められる特有さの情報から求めた各用語及び用語対の重みとを用いて求め、一定値以上の文書共通度を有する文書を着目クラスターに一時的に帰属させるステップと、
- (i) 上記 (f) から (h) までのステップを着目クラスターに一時的に帰属する文書数が増えなくなるまで繰返すステップと、
- (j) 上記 (d) から (i) までのステップを、どのクラスターに対しても文書共通度が一定値以下の文書の数が 0 になるか、もしくはその数が一定値以下で前回の繰返しの時と同じになるまで繰返すステップと、
- (k) 各文書の各クラスターに対する文書共通度をもとに、各文書の帰属するクラスターを決定するステップと、
- (l) 冗長なクラスターの有無をチェックし、冗長なクラスターが存在すれば除去したうえで、各文書の帰属するクラスターを改めて決定するステップ。

【請求項 6】

前記出現する用語の種類数が M で与えられ、 R 個の文書からなる文書集合 D において、 r 番目の文書を D_r 、 D_r の文書セグメント数を Y_r 、 D_r の y 番目の文書セグメントベクトルを d_{ry} (d_{ry1}, \dots, d_{ryM})^T とすると、ここで、 T はベクトルの転置を表す、文書 D_r の前記共起行列 S^r は、

【数 1】

$$S^r = \sum_{y=1}^Y d_{ry} d_{ry}^T \quad \dots \dots (1)$$

で与えられることを特徴とする請求項 5 に記載の方法。

【請求項 7】

文書集合 D の文書頻度行列の各成分は、文書集合 D 中の各文書の共起行列の対応する成分がゼロでない文書数であることを特徴とする請求項 6 に記載の方法。

【請求項 8】

文書集合 D の共通共起行列は、m n 成分が以下のように決定される行列 T

【数 2】

$$T_{mn} = \prod_{r=1}^R S^r_{mn} \quad S^r_{mn} > 0 \quad \dots \dots (2)$$

をもとに、m n 成分が

$$T^A_{mn} = T_{mn}, \quad U_{mn} \geq A \text{ の時}$$

$$T^A_{mn} = 0 \quad \text{それ以外のとき}$$

によって決定される行列 T^A により、もしくは m n 成分が

$$Q^A_{mn} = \log(T^A_{mn}) \quad T^A_{mn} > 1 \text{ の時}$$

$$Q^A_{mn} = 0 \quad \text{それ以外のとき}$$

によって決定される行列 Q^A により与えられることを特徴とする請求項 6、7 に記載の方法。

【請求項 9】

共起行列を S^P とする文書 P の文書集合 D に対する文書共通度は、 z_{mm} 、 z_{mn} をそれぞれ用語 m、用語対 m、n に対する重みとして

【数 3】

$$com_1(D, P; Q^A) = \frac{\sum_{m=1}^M z_{mm} Q^A_{mm} S^P_{mm}}{\sqrt{\sum_{m=1}^M z_{mm} (Q^A_{mm})^2} \sqrt{\sum_{m=1}^M z_{mm} (S^P_{mm})^2}} \quad \dots \dots (3)$$

もしくは、

【数 4】

$$com_q(D, P; Q^A) = \frac{\sum_{m=1}^M \sum_{n=1}^M z_{mn} Q_{mn}^A S_{mn}^P}{\sqrt{\sum_{m=1}^M \sum_{n=1}^M z_{mn} (Q_{mn}^A)^2} \sqrt{\sum_{m=1}^M \sum_{n=1}^M z_{mn} (S_{mn}^P)^2}}$$

... (4)

もしくは数式 (3)、数式 (4) において行列 Q^A の代わりに行列 T^A を用いた式により与えられることを特徴とする請求項 6、7、8 に記載の方法。

【請求項 10】

クラスターの種となる文書の抽出と初期クラスターの作成は以下の (a) から (d) のステップを経て行われることを特徴とする請求項 6、7、8、9 に記載の方法。

(a) その時点で求められているどのクラスターに対しても文書共通度が一定値以下の文書の集合から求められる共通共起行列を用いて、前記文書集合中の各文書の文書集合共通度、もしくは情報共通量を求めるステップと、

(b) (a) で求められた文書集合共通度、もしくは情報共通量の大きい一定個の文書をクラスターの種の候補として抽出するステップと、

(c) クラスターの種の各候補について、前記文書集合中の各文書との類似度を求め、類似度が一定値以上となる文書を近隣文書として求めるステップと、

(d) クラスターの種の各候補の中から近隣文書の最も多い候補を選択してクラスターの種の文書とし、その近隣文書を初期クラスターとするステップ。

【請求項 11】

共起行列を S^P とする文書 P の文書集合 D に対する情報共通量は

【数 5】

$$comInfo(D, P) = \sum_{m=1}^M \sum_{n=1}^M z_{mn} Q_{mn}^A S_{mn}^P$$

... (5)

もしくは数式 (5) において行列 Q^A の代わりに行列 T^A を用いた式により与えられることを特徴とする請求項 10 に記載の方法。

【請求項 12】

各用語及び用語対の着目クラスターに対する特有さの情報と重みの決定は、以下の (a) から (d) のステップを経て行われることを特徴とする請求項 6、7、8、9、10、11 に記載の方法。

(a) 入力文書全体から求められる文書頻度行列の各成分の、その時点で着目クラスターに一時的に帰属する文書集合から求められる文書頻度行列の対応する成分に対する比を、対角成分の場合には用語文書頻度比として、非対角成分の場合には用語対文書頻度比として求めるステップと、

(b) その時点で着目クラスターに一時的に帰属する文書集合において、最も頻度の高い一定個の用語の中で用語文書頻度比が小さい一定個の用語を選択し、それらの用語文書頻度比の平均を平均文書頻度比として求める、もしくは、最も頻度の高い一定個の用語や用語対の中で用語もしくは用語対文書頻度比が小さい一定個の用語や用語対を選択し、それらの用語もしくは用語対文書頻度比の平均を平均文書頻度比として求めるステップと、

(c) 平均文書頻度比を各用語または用語対の用語文書頻度比または用語対文書頻度比で除した値を各用語または用語対の特有さの情報として求めるステップと

(d) 特有さの情報を変数とする関数によって用語または用語対の重みを決定するステップ。

【請求項 13】

文書入力部、文書前処理部、文書情報処理部、及び出力部を有する装置において、一つまたは複数の文書セグメントを持つ複数の文書から成る入力文書集合に対して、文書間の関係に基づき文書をグループ化するために、

(a) 前記文書集合と共通性の高い文書を前記文書集合から選択しクラスターとするクラスター選択部、

(b) 前記選択された文書と類似性の高い全ての文書を前記クラスターに追加する追加部、

(c) 前記クラスターに存在する文書と共通性の高い全ての文書を前記クラスターに繰返し追加するクラスター成長部、

(d) 既に求められたクラスターと共通性の高い文書を除外した文書集合を前記

文書集合として前記 (a)、(b) 及び (c) のステップを繰返す繰返し部、
(e) 前記クラスターが重複する場合に、冗長なクラスターを除去する除去部、
を動作させる為のプログラム。

【請求項 14】

文書入力部、文書前処理部、文書情報処理部、及び出力部を有する装置において、以下の (a) から (d) の手段を有する、一つまたは複数の文書セグメントを持つ複数の文書から成る文書集合に対して、文書間の関係に基づき文書をグループ化するクラスタリング装置、

(a) 前記文書集合と共通性の高い文書を前記文書集合から選択しクラスターとする手段、

(b) 前記選択された文書と類似性の高い全ての文書を前記クラスターに追加する手段、

(c) 前記クラスターに存在する文書と共通性の高い全ての文書を前記クラスターに繰返し追加する手段、

(d) 既に求められたクラスターと共通性の高い文書を除外した文書集合を前記文書集合として前記 (a)、(b) 及び (c) のステップを繰返す手段、

(e) 前記クラスターが重複する場合に、冗長なクラスターを除去する手段。

【発明の詳細な説明】

【0001】

【産業上の利用分野】

本発明は文書のクラスタリングをはじめとする自然言語処理に関するものであり、前記処理の高性能化を図ることによって文書からの情報の抽出を容易にするものである。

【0002】

【従来の技術】

文書クラスタリングは入力された文書群を文書の内容によって幾つかのグループに分割する技術である。クラスタリング技術は古くから研究されてきており、これまでに考案された方法については C. D. Manning と H. Schütze によって著された Foundations of Statistical

Natural Language Processing (The MIT Press、1999年) に体系的に紹介されている。先ず、クラスタリングは、各文書が各クラスターに帰属する確率を求めるソフトクラスタリング、各クラスターに帰属するか否かを求めるハードクラスタリングに大別される。後者については、さらに、階層的な手法と非階層的な手法とに分類される。階層的な手法は、さらにボトムアップのアプローチとトップダウンのアプローチに分けられる。前者では、初期状態として各文書をクラスターの核とし、最も近いクラスター同士をマージするという処理を繰り返す。これにより文書集合は木構造で表現されるようになる。クラスター間の近さの尺度、即ち類似度を図る方法として単一リンク法、完全リンク法、グループ平均法が知られている。これらは何れも2文書間の類似度をもとに算出されるものである。後者では、全文書が1つのクラスターに属するという状態から出発し、例えばひとつのクラスター中のあらゆる文書対の中で最も低い類似度が閾値以下の場合、そのクラスターを分割するという処理を繰り返す。

【0003】

非階層的な手法では、予め指定された数のクラスターが何らかの基準を満たすように作成される。よく知られている方法の例を以下に示す。

ステップ1：指定されたクラスター数の文書をランダムに選択して各クラスターの中心とするステップ、

ステップ2：各文書について各クラスター中心との近さを求め、各文書を最も近いクラスターに帰属させるステップ、

ステップ3：各クラスターに帰属する文書ベクトルの平均により各クラスターの中心を求めるステップ、

ステップ4：ステップ2) の処理を実行し、各文書の帰属するクラスターに変化がなければ終了、そうでなければステップ3へいくステップ、からなる方法である。

【0004】

【発明が解消しようとする課題】

従来の文書クラスタリング技術は2つの大きな問題を抱えている。ひとつは求め

られるクラスターの数の問題である。文書クラスタリングでは求められるクラスターの数は入力された文書集合中の文書が述べている話題の数と同じでなければならない。前述のようにボトムアップの階層的なクラスタリング処理では、各クラスターは一つの文書から成る状態から出発し、最も近いクラスター同士をマージする処理を繰り返して最後は全文書が一つのクラスターに属することになる。従って、話題の数と同じ数のクラスターを得るにはクラスターのマージを打ち切ることが必要となる。これは、クラスターのマージ処理において類似度が閾値以下のクラスター対についてはマージを行わないようにすることにより実現可能である。しかしながら、実際には上述の閾値を如何に決めるかが難しい問題となっている。閾値が不適切であれば正しい数のクラスターは得られない。同様にトップダウンのクラスタリング処理では、ひとつのクラスター中のあらゆる文書対の中で最も低い類似度が閾値以上の場合にはクラスターの分割は行わないようにすれば、原理的には話題の数と同じ数のクラスターが得られる筈である。

【0005】

しかしながら、この場合にも閾値をどのように決めるかが難しい問題となっている。また、非階層的なクラスタリングでは、ユーザーは与えられた文書集合を何個のクラスターに分割するかを事前に入力することを求められる。しかしながら、クラスター数の情報は入力文書集合の事前の知識なしには正確に与えることは不可能である。このように入力文書集合から正しい数のクラスターを得ることは難しい問題となっている。非階層的なクラスタリングでクラスター数を正しく推測する試みもなされており、性能は向上しているが完璧ではない (X. Liu, Y. Gong, W. Xu and S. Zhu. Document Clustering with Cluster Refinement and Model Selection Capabilities. In Proceedings of the 25th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, pp. 191-198. Tampere, Finland, August, 2002.)。

【0006】

2番目の問題はクラスタリングにおける精度の問題である。これは、ひとつのクラスターに属する文書が同じ話題を述べているかどうかの問題である。クラスタリング処理では、通常は文書は、各成分が文書中に現れる各用語の有無、もしくは出現頻度で与えられるベクトルで表現され、2つのクラスター間の類似度は、異なるクラスターに属する文書対の余弦類似度をもとに、ある文書とクラスター間の類似度はその文書のベクトルとそのクラスターに属する文書の平均ベクトルとの間の距離（例えばユークリッド距離）をもとに求められていた。従来のクラスタリング処理では、上記の余弦類似度やユークリッド距離を求めるときに、どのような用語がそのクラスターに重要なのかを検証することなく、各文書で得られたベクトルをそのまま用いることが多かった。そのため、各クラスターに本質的でない情報を表す用語や用語対の存在がクラスタリングの精度に影響を及ぼした。

【0007】

【課題を解決するための手段】

ところで、文書クラスタリングは各文書が記述する話題によって文書をグループ化するものであるから、ひとつのクラスターに属する文書（クラスター文書集合と呼ぶ）は同じ話題について述べている筈である。従って、クラスター文書集合は何らかの共通性を有する筈である。また、各話題にはその話題には多く出現し、他の話題にはあまり出現しない話題特有の用語や用語対が存在する筈である。従って、各クラスターには用語や用語対の出現傾向に違いが存在する筈である。このようなことから、本発明ではクラスタリングの精度を高めるために、クラスタリングの過程で、

- A)：着目クラスター文書集合の共通情報を抽出し、この共通情報との近さにより各文書の着目クラスターに対する近さ（文書共通度）を求める。
- B)：着目クラスターに特有でない用語や用語対を検出し、上記文書共通度の算出においてはそれらの影響を排除する。

という手段を導入する。

【0008】

従来の階層的な処理では、クラスターのマージや分割が頻繁に繰り返される。また、従来の非階層的な処理では、クラスターのメンバーがやはり頻繁に入れ替わる。このような状況では各クラスターの共通情報、クラスターに非特有な用語や用語対の検出には無理がある。そこで、本発明ではクラスタリングの全体の手順として以下を採用する。

【0009】

ステップ1：最初の繰り返しにおいては全文書から、2回目以降の繰り返しにおいては、その時点で存在するどのクラスターに対しても文書共通度が一定値以下の文書の中から、クラスターの種の候補となる文書を検出する。

ステップ2：先ず、クラスターの種の各候補文書について、全文書との類似度を求め一定値以上の類似度を有する文書を抽出する。その文書数が最も大きくなる候補文書をクラスターの種とし、その文書集合によりクラスターを形成する。

ステップ3：その時点でのクラスター文書集合と各文書との間で文書共通度を求め、一定値以上の文書共通度を有する文書をそのクラスターに一時的に帰属させることによりクラスターを成長させる。クラスターに一時的に属する文書数が一定になればステップ4へ移行する。そうでなければ本ステップを繰り返す。

ステップ4：終了条件を満たせばステップ5へ移行する。そうでなければステップ1に戻って続行する。

ステップ5：各文書について各クラスターへの文書共通度を求め、文書共通度の最も高くなるクラスターに帰属させる。

ステップ6：1つの話題に2つ以上のクラスターが対応していないかどうかを検出する。そのようなクラスターがあれば冗長なクラスターとして削除し、各文書の帰属するクラスターを求めなおす。

【0010】

上記のクラスタリング手順において、前記A)の共通情報を用いる文書共通度の算出、及びB)の着目クラスターに特有でない用語や用語対の検出はステップ3と5において行われることになる。A)については、その時点での着目クラスターに一時的に帰属している文書から共通情報を抽出することになる。共通情報の抽出と利用は特願2002-326157で述べられている方法を援用すること

ができる。基本的な考え方は次の通りである。いま、着目クラスターがR個の文書から成るものとし、各文書から一つずつ文を取り出してR個の文からなる文の組を作ったとする。このような文の組は各文書の文の数の積通り存在することになる。ここでは、着目する文の組において、R個の文のうちのA個の文に現れる用語を共通用語、共通用語で構成された文を共通文と呼ぶこととする。ここで、全ての文の組で共通文を作り、共通文の集合を構成したとする。このような共通文の集合は着目クラスターの共通の話題の内容を示すものと考えられる。従って、各文書と共通文集合との間で何らかの手段で類似度を求めることができれば、それは各文書の着目クラスターの共通話題への近さを表わすことになる。

【0011】

またB)の着目クラスターに特有でない用語や用語対の検出は次のような考え方で行うものである。種文書が話題iの着目クラスターの成長の過程を考える。話題iを述べている文書数は、文書集合全体には c_0 個、着目クラスターの文書集合には c 個存在したとする。また、用語mを含む文書数は、文書集合全体では U_{mm}^0 個、着目クラスターの文書集合では U_{mm} 個存在したとする。用語mが話題iの特有用語の時には用語mは話題iの文書のみに見えるので、

【0012】

【数6】

$$U_{mm}^0 / U_{mm} \approx c_0 / c$$

【0013】

となり、非特有の時には話題i以外の話題の文書にも現れるので、

【0014】

【数7】

$$U_{mm}^0 / U_{mm} > c_0 / c$$

【0015】

となる筈である。従って、 c_0 / c を適当な方法で求めることができれば用語mが話題iに特有か否かを判断することができる。 U_{mm}^0 / U_{mm} を用語mの文

書頻度比と呼ぶことにすると、本発明では、着目クラスターの文書集合において最も頻度の高い一定個の用語のうち、文書頻度比の値の小さな一定個は話題 i の特有用語とみなし、これらの用語の文書頻度比の平均 c' を c_0/c の推測値とした。結局、 α をパラメータとして

【0016】

【数8】

$$U_{mm}^0/U_{mm} > \alpha c'$$

【0017】

を満たす用語 m は話題 i には特有用語ではないと判断できる。

同様に、用語 m 、 n を含む文書数は、文書集合全体では U_{mn}^0 個、着目クラスターの文書集合では U_{mn} 個存在したとして、

【0018】

【数9】

$$U_{mn}^0/U_{mn} > \alpha c'$$

【0019】

を満たす用語対 m 、 n は話題 i には特有用語対ではないと判断できる。

文書共通度は着目クラスターに本質的でない用語や用語対の影響を受けにくくするためには、話題 i には特有ではないと判断された用語、用語対は、各文書と着目クラスターの文書集合との文書共通度の算出に用いないようにすればよい。もしくは、

【0020】

【数10】

$$c/(U_{mm}^0/U_{mm})$$

【0021】

【数11】

$$c/(U_{mn}^0/U_{mn})$$

【0022】

をそれぞれ用語 m 、用語対 m 、 n の重みとして用いて文書共通度の算出を行ってもよい。このようにすることにより、話題 i を述べた文書に対して文書共通度は

大きな値をとるようになる。その結果クラスタリングの精度の向上が期待できる。

【0023】

上記のクラスタリングの全体手順においては、先ずクラスターの種となる文書をひとつ取り出し、ついでその種と同じ話題を記述する文書を検出して種を成長させるという処理を繰り返し行うこととなる。従って、種の文書の数が入力文書における話題の数と過不足なく一致すれば正しい数のクラスターが得られることになる。たとえステップ1において同じ話題に対して2つの種文書が検出されたにしても、ステップ6で冗長なクラスターを検出して除去するので正しい数のクラスターが得られる。また、ステップ1においてある話題に対して種文書が検出されない時にはクラスターの数不足することになる。このような事態は、クラスタリングの精度が低く、検出されるべき話題の文書がすでに存在するクラスターとの文書類似度が高くなってしまった時に起きる。言い換えれば、このような事態は、クラスタリングの精度が低いために、ひとつのクラスターに本来有すべきでない話題を有する文書が混入することによって引き起こされる。しかしながら、本発明では上記A)、B)の手段によってクラスタリングの精度を高めているので異なる話題の文書の混入の可能性は低く、求められるクラスターの数が少なくなるという事態は起こりにくい。本発明では正しい数のクラスターが得られる公算が非常に大きい。

【0024】

【実施例】

図1は、本発明の概要を示すブロック図である。110は文書入力ブロック、120は文書前処理ブロック、130は文書情報処理ブロック、140は出力ブロックを示す。文書入力ブロック110には、処理したい文書集合が入力される。文書前処理ブロック120では、入力された文書の用語検出、形態素解析、文書セグメント区分け等が行われる。文書セグメントについて説明する。文書セグメントは文書を構成する要素であり、その最も基本的な単位は文である。英文の場合、文はピリオドで終わり、その後ろにスペースが続くので文の切出しは容易に行うことができる。その他の文書セグメントへの区分け法としては、ひとつの文

が複文からなる場合主節と従属節に分けておく方法、用語の数がほぼ同じになるように複数の文をまとめて文書セグメントとする方法、文書の先頭から含まれる用語の数が同じになるように文とは関係なく分けする方法などがある。文書情報処理ブロック 130 は以下に詳細に説明するが、情報処理を行い、種の文書の検出、全文書と着目クラスターとの文書集合共通度の算出、各クラスターの特有でない用語、用語対の検出などクラスタリングに直接関わる処理を行う。出力ブロック 140 は文書情報処理ブロック 130 で得られた結果を、ディスプレイ等の出力装置に出力する。

【0025】

図2は与えられた文書集合に対して、クラスタリングを行う本発明の実施例を示す。この発明の方法は、汎用コンピュータ上でこの発明を組み込んだプログラムを走らせることによって実施することができる。図2は、そのようなプログラムを走らせている状態でのコンピュータのフローチャートである。ブロック21は文書集合入力、ブロック22は文書前処理、ブロック23は全体文書情報抽出処理、ブロック24はクラスターの種文書及び初期クラスターの決定、ブロック25はクラスターの成長処理、ブロック26は残存文書検出、ブロック27は終了条件チェック、ブロック28は残存文書の文書情報抽出処理、ブロック29は帰属クラスター決定、ブロック30は冗長クラスターの検出・除去、である。以下、英文文書を例に実施例を説明する。

【0026】

まず、文書集合入力21において対象となる文書集合が入力される。文書前処理22においては各入力文書に対して、用語検出、形態素解析、文書セグメント分け、文書セグメントベクトル作成などの前処理が行われる。用語検出としては、各入力文書から単語、数式、記号系列などを検出する。ここでは、単語や記号系列などを総称して全て用語と呼ぶ。英文の場合、用語同士を分けて書く正書法が確立しているので用語の検出は容易である。次に、形態素解析では、各入力文書に対して用語の品詞付けなどの形態素解析を行う。文書セグメント分けにおいて各入力文書に対して文書セグメントへの分けを行う。文書セグメントベクトル作成では、まず文書全体に出現する用語から作成すべきベクトルの次元数

および各次元と各用語との対応を決定する。この際に出現する全ての用語の種類にベクトルの成分を対応させなければならないということはなく、品詞付け処理の結果を用い、例えば名詞と動詞と判定された用語のみを用いてベクトルを作成するようにしてもよい。次いで、各文書セグメントに出現する用語に対応する成分のみが値1、他は0となるような文書セグメントベクトルを作成する。

【0027】

全体文書情報抽出処理23では、後段のクラスタリング処理に用いるデータを各文書と入力文書集合全体から求める。求めるデータは各文書の共起行列、前記共通文集合の共起行列（共通共起行列）、入力文書集合全体の文書頻度行列である。各文書の共起行列は用語の出現頻度、用語間の共起頻度を反映する行列である。文を文書セグメントとした場合について説明を続ける。ここでは、現れる用語集合が $\{w_1, \dots, w_M\}$ で与えられ、 R 個の文書から成る入力文書集合 D を考える。さらに、 r 番目の文書を D_r とすると、 D_r は Y_r 個の文からなるものとし、 y 番目の文及びその文ベクトルを D_{ry} 、 $d_{ry} = (d_{ry1}, \dots, d_{ryM})^T$ とする。ここで、 T はベクトルの転置を表す。 d_{ry} は2値ベクトルであり、 d_{rym} は m 番目の用語の有無を表す。文書の D_r の共起行列を S^r とすると、これは

【0028】

【数12】

$$S^r = \sum_{y=1}^{Y_r} d_{ry} d_{ry}^T \quad \dots (1)$$

【0029】

$S^r_{mn} = \sum_{y=1}^{Y_r} d_{rym} d_{ryn}$ で与えられる。式(1)から分かるように、 S^r の mn 成分はに

【0030】

【数13】

$$T_{mn} = \prod_{r=1}^R S^r_{mn} \quad S^r_{mn} > 0$$

【0031】

さらに、各用語、用語共起の文書頻度を格納した行列 U^0 を求める。 U^0_{mm} 、 U^0_{mn} はそれぞれ用語 m の出現する文書数、用語 m 、 n の共起する文書数となる。このような行列 T 、 U^0 を用いて共通共起行列 T^A を求める。共通共起行列 T^A の mn 成分は以下のように決められる。

$$T^A_{mn} = T_{mn}, \quad U^0_{mn} \geq A \text{ のとき}$$

$$T^A_{mn} = 0 \quad \text{それ以外}$$

A はパラメータであり、実験的に決められる。

【0032】

また、 mn 成分が以下のように与えられる行列 Q^A を

$$Q^A_{mn} = \log(T^A_{mn}) \quad T^A_{mn} > 1 \text{ のとき}$$

$$Q^A_{mn} = 0 \quad \text{それ以外}$$

により定義し、共通共起行列として用いてもよい。

【0033】

クラスターの種文書及び初期クラスターの決定 24 では、前記ステップ 1 及び 2 に対応する処理を行う。ここで、その時点で存在するどのクラスターに対しても文書共通度が一定値以下の文書の集合を D' とする。文書集合 D' はその時点で存在するどのクラスターにも属さない公算の大きい文書の集合である。共通共起行列 T^A 、 Q^A 、文書頻度行列 U は最初の繰り返しにおいては全文書 D に基づいて計算され、2 回目以降の繰り返しにおいては、文書集合 D' に基づいて計算される。一方、クラスターの種となる文書はどの話題の文書が選ばれようと、その話題の中では中心的な文書であることが望ましい。本発明では、 D' の中で最も優勢な話題の文書群において中心的な文書は、 D' との共通度も高いであろうとの仮定のもとに、文書集合 D' 中の文書と文書集合 D' との文書共通度を求め、文書共通度の高い文書をクラスターの種の候補として選択する。任意の文書を P 、その共起行列を S^P とするとき、文書 P と D' との文書共通度としては例えば以下を用いることができる。

【0034】

【数 14】

$$com_q(D', P, Q^A) = \frac{\sum_{m=1}^M \sum_{n=1}^M Q_{mn}^A S_{mn}^P}{\sqrt{\sum_{m=1}^M \sum_{n=1}^M (Q_{mn}^A)^2} \sqrt{\sum_{m=1}^M \sum_{n=1}^M (S_{mn}^P)^2}}$$

... (2)

【0035】

式(2)において行列 Q^A の代わりに行列 T^A を用いることもできる。また、式(2)において、複数の話題に共通する用語の影響を軽減するために、共起行列、共通共起行列の対角成分は用いないようにしてもよい。

【0036】

クラスターの種文書の候補は、 D' 中の全ての文書に対して式(2)により文書共通度を求め、文書共通度の高い一定個の文書を選択することにより求められる。次にクラスターの種文書の決定について説明する。まず各候補文書について D' 中の全ての文書との類似度を求める。次いで各候補文書について各候補文書との類似度が一定値より大きい文書を各候補文書の近隣文書として求める。近隣文書数が最も大きい文書を候補文書の中から一文書選択することにより、クラスターの種文書が決定される。また、初期クラスターはその種文書の近隣文書で与えられる。

【0037】

クラスターの成長処理25では、クラスターの初期クラスターと共通度の高い文書を吸収することによりクラスターを成長させる。図3はそのようなクラスターを成長させる処理のブロック図である。31は文書頻度行列作成、32は共通共起行列作成、33は用語・用語対の特有度算出、34は文書共通度算出、35はクラスターメンバー決定、36は終了条件チェックである。

【0038】

文書頻度行列作成31、共通共起行列作成32では、その時点で一時的に着目クラスターのメンバーとなっている文書集合に対して、図2のブロック23における文書頻度行列作成処理、共通共起行列作成と同等の処理を行う。31で求め

られた文書頻度行列を U により表す。32で求められた共通共起行列を TA もしくは QA により表す。用語・用語対の特有度算出33では、各用語・用語対の特有度を決定し、重みを決定する。まず、前述のように、 U^0_{mm}/U_{mm} を用語 m の文書頻度比として求め、 U から求められる最も頻度の高い一定個の用語のうち、文書頻度比の値の小さな一定個を着目クラスターの特有用語として選択する。次いで、これらの用語の文書頻度比の平均を平均文書頻度比 c' として求め、用語 m の特有度 v_{mm} 、用語対 m, n の特有度 v_{mn} を以下により決定する。

【0039】

【数15】

$$v_{mm} = c' / (U^0_{mm} / U_{mm})$$

【0040】

【数16】

$$v_{mn} = c' / (U^0_{mn} / U_{mn})$$

【0041】

あるいは、特有用語に限定することなく、特有用語対と特有用語の両方を用いて平均文書頻度比を求めるようにしてもよい。この場合には U^0_{mn}/U_{mn} を m と n が等しくないときは用語対 m, n の文書頻度比として、 m と n が等しいときは用語 m の文書頻度比として求め、 U から求められる最も頻度の高い一定個の用語、用語対のうち、文書頻度比の値の小さな一定個を着目クラスターの特有用語または用語対として選択する。次いで、これらの用語、用語対の文書頻度比の平均を平均文書頻度比 c' として求めるようにする。

【0042】

用語 m 、用語対 m, n の重みをそれぞれ z_{mm} 、 z_{mn} とする。これらは、重み決定関数 $f(x)$ を用いて、以下のように決定する。

【0043】

【数17】

$$z_{mm} = f(v_{mm})$$

【0044】

【数18】

$$z_{mn} = f(v_{mn})$$

【0045】

$f(x)$ の決め方は任意であるが、 $f(x) = x$ 、 $f(x) = x^2$ のようにするのがひとつの方法である。あるいは、 x が一定値よりも大きいときは $f(x) = 1$ 、そうでないときは $f(x) = 0$ としてもよい。

【0046】

文書共通度算出34では、全入力文書に対して着目クラスターとの文書共通度を算出する。任意の文書を P 、その共起行列を S^P とするとき、文書 P の文書共通度は、

【0047】

【数19】

$$com_1(D, P; Q^A) = \frac{\sum_{m=1}^M z_{mm} Q_{mn}^A S_{mn}^P}{\sqrt{\sum_{m=1}^M z_{mm} (Q_{mn}^A)^2} \sqrt{\sum_{m=1}^M z_{mm} (S_{mn}^P)^2}}$$

【0048】

もしくは、

【0049】

【数20】

$$com_q(D, P; Q^A) = \frac{\sum_{m=1}^M \sum_{n=1}^M z_{mn} Q_{mn}^A S_{mn}^P}{\sqrt{\sum_{m=1}^M \sum_{n=1}^M z_{mn} (Q_{mn}^A)^2} \sqrt{\sum_{m=1}^M \sum_{n=1}^M z_{mn} (S_{mn}^P)^2}}$$

【0050】

により求めることができる。上式において行列 Q^A の代わりに行列 T^A を用いることもできる。

【0051】

クラスターメンバー決定35では、ブロック34で求められた各文書の着目クラスターに対する文書共通度を用いて、一定値以上の文書共通度を有する文書を求め、着目クラスターの一時的なメンバーとする。

【0052】

終了条件チェック36では、着目クラスターの成長処理を終了するか否かをチ

ェックする。まず、初回の繰り返し、即ち初めて36に到達したときには無条件に31に戻って処理を繰り返す。2回目以降の繰り返しの時には、上記35で求められた着目クラスターのメンバーの文書数をカウントし、それが前回の繰り返しのメンバー数と等しくない場合には31に戻って処理を繰り返す。等しければ、全入力文書と着目クラスターとの文書共通度を保持し、着目クラスターの成長処理を終了する。

【0053】

図2に戻って説明を続ける。残存文書検出26では、その時点で存在する全てのクラスターに対する各文書の文書共通度をもとに、どのクラスターに対しても文書共通度が一定値以下の文書を残存文書として抽出する。

【0054】

終了条件チェック27では、残存文書数をもとにクラスターの種の検出から成長に到る一連の処理を終了するか否かのチェックを行う。例えば、残存文書数が0、もしくは残存文書数が一定値以下でかつ前回の繰り返し時の残存文書数と等しい場合にはブロック29に移行するようにする。もし、このような条件が満たされなければ、ブロック28に移行し、残存文書集合に対してブロック23と同等な処理を行う。

【0055】

帰属クラスター決定29では、各文書が最終的に求められたクラスターのどれに帰属するかを決定する。これは、図3において各文書について求められた各クラスターに対する文書共通度の情報を用い、各文書は文書共通度が最も高くなるクラスターに帰属させることで実行できる。

【0056】

冗長クラスターの検出・除去30では、冗長なクラスターが存在するか否かをチェックし、存在する場合には除去する。冗長なクラスターは、ひとつの話題に対して2つ以上のクラスターが求められたときに発生する。そのような時、その話題を記述した文書は2つ以上のクラスターに対して大きな文書共通度を有するようになり、2つ以上のクラスターが重複する格好になる。冗長なクラスター検出のためには、まず、求められた全てのクラスターと全ての文書との文書共通度を

求め、次いで文書共通度が着目クラスターに対しては一定値よりも大きく、他のどのクラスターに対しても一定値よりも小さくなる文書数を求める。そのような文書数は着目クラスターが他のクラスターと重複しない場合には、その着目クラスターと一定値以上の文書共通度を有する文書数と等しくなる。一方、着目クラスターが他のクラスターと重複する場合には、他のクラスターと重複していない文書数すなわち当該着目クラスターにのみ属する文書数となる。このような文書数を各クラスターの重要度と定義する。重要度は、図4 (a) の場合「クラスター1」及び「クラスター2」に属する文書数である。重要度は、着目クラスターが他のクラスターと一部重複して存在する場合には、他のクラスターと重複していない文書数となる。すなわち、「クラスター1」に対しては、図4 (b) の「c」で示される部分に含まれる文書数を意味する。「クラスター2」に対しては、図4 (b) の「d」で示される部分に含まれる文書数を意味する。ひとつのクラスターの重要度が一定値よりも小さい場合は、そのクラスターに属する文書の数がいくら多くとも冗長なクラスターとみなし除去する。そのようなクラスターが複数存在すれば、クラスター重要度が最も小さいクラスターを先ず除去する。その後残されたクラスターについてクラスター重要度の算出を行い、クラスター重要度が最も小さいクラスターを除去する。このような処理を冗長なクラスターが存在しなくなるまで繰り返す。クラスターの削除を行った場合は、除去されたクラスターに属していた文書の帰属クラスターの決定を改めて行う。

【0057】

【発明の効果】

ここで本発明の効果を説明する為に図2、3の実施例に沿った実験結果を示す。用いたコーパスはTDT2である。TDT2は1998年の1月から6月の間の100個のイベントに関するニュースストーリーの集合であり、6個のニュースソースから採取されている。本報告では同じくTDT2を用いて行われたLiuらの非階層型のクラスタリング (X. Liu, Y. Gong, W. Xu and S. Zhu. Document Clustering with Cluster Refinement and Model Selection Capabilities. In Proceedings of the 25

th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, pp. 191-198. Tampere, Finland, August, 2002.) の結果と比較するため、L i uらが行ったようにABC、CNN、VOAから採取された15イベントに関するニュースストーリーの集合を実験対象とした。表1にそれらの詳細を示す。

【0058】

【表1】

イベント ID	イベントの内容	文書数			
		ABC	CNN	VOA	Total
01	Asian Economic Crisis	27	90	289	406
02	Monica Lewinsky Case	102	497	96	695
13	1998 Winter Olympic	21	81	108	210
15	Current Conflict with Iraq	77	438	345	860
18	Bombing AL Clinic	9	73	5	87
23	Violence in Algeria	1	1	60	62
32	Sgt. Gene McKinney	6	91	3	100
39	India Parliamentary Election	1	1	29	31
44	National Tobacco Settlement	26	163	17	206
48	Jonesboro shooting	13	73	15	101
70	India, A Nuclear Power?	24	98	129	251
71	Israeli-Palestinian Talks	5	62	48	115
76	Anti-Suharto Violence	13	55	114	182
77	Unabomber	9	66	6	81
86	GM Strike	14	83	24	121

【0059】

表2に実験に用いられた14種類のデータセットとそれに対する提案手法とL i uらの手法のクラスタリングの精度を示す。L i uらの手法の結果はL i uらの論文より再掲したものである。本発明では、ある文書が属するイベントとその文書が帰属するクラスターの種となった文書のイベントが一致するときクラスタリングの結果は正しいとされる。また、全てのクラスターに対して文書共通度が

0の文書は誤りとする。精度は正しくクラスタリングされた文書数の全文書数に対する比により求める。L i uらの方法は、混合ガウス分布モデル (G a u s s i a n M i x t u r e M o d e l) に基づき非階層形のクラスタリングを行った後、各クラスターの特有用語を求め、特有用語の投票によって結果を修正している。表2において、テストデータのABC-01-02-15とあるのは、ABCより採取されたイベントIDが01、02、15に属する文書であることを意味している。表2から、精度の高いデータセットの数は、L i uらの方法よりも本発明の方が多く、本発明が優ることが分かる。

【0060】

【表2】

番号	データセット	Liuらの方法	本発明
1	ABC-01-02-15	1.0000	0.9806
2	ABC-02-15-44	0.9902	0.9805
3	ABC-01-13-44-70	1.0000	1.0000
4	ABC-01-44-48-70	1.0000	1.0000
5	CNN-01-02-15	0.9756	0.9932
6	CNN-02-15-44	0.9964	0.9964
7	VOA-01-02-15	0.9896	0.9986
8	VOA-01-13-76	0.9583	0.8943
9	VOA-01-23-70-76	0.9453	0.9206
10	VOA-12-39-48-71	0.9898	1.0000
11	VOA-44-48-70-71-76-77-86	0.8527	1.0000
12	ABC+CNN-01-13-18-32-48-70-71-77-86	0.9704	0.9917
13	CNN+VOA-01-13-48-70-71-76-77-86	0.9262	0.9500
14	ABC+CNN+VOA-44-48-70-71-76-77-86	0.9938	1.0000

【0061】

また、求められるクラスターの個数も表2のデータに対して全て正しく求められている。

また、L i uらの論文で挙げられている12種類のデータに対してもクラスター数は正しく求められている。一方、L i uらの方法では12種類のうち3種類のデータに対して正しく求められていない。表3にL i uらの方法及び本発明の結果を示す。

【表3】

テストデータ	正しい クラスタ数	Liu等の テスト結果	本願発明の テスト結果
ABC-01-03	2	2	2
ABC-01-02-15	3	3	3
ABC-02-48-70	3	2	3
ABC-44-70-01-13	4	4	4
ABC-44-48-70-76	4	4	4
CNN-01-02-15	3	4	3
CNN-01-02-13-15-18	5	5	5
CNN-44-48-70-71-76-77	6	5	6
VOA-01-02-15	3	3	3
VOA-01-13-76	3	3	3
VOA-01-23-70-76	4	4	4
VOA-12-39-48-71	4	4	4

【0062】

このように本発明によれば、入力文書集合から正しい数のクラスターを抽出すること、及び各文書が帰属すべきクラスターを精度よく決定することができ、ユーザの情報獲得の効率性が高められる。

【図面の簡単な説明】

【図1】

本発明の概略を示すブロック図である。

【図2】

文書集合が入力された段階からクラスターと各文書が帰属するクラスターが決定されるまでの手順を示す図である。

【図 3】

一つのクラスターについて初期クラスターの段階からの成長の手順を示す図である。

【図 4】

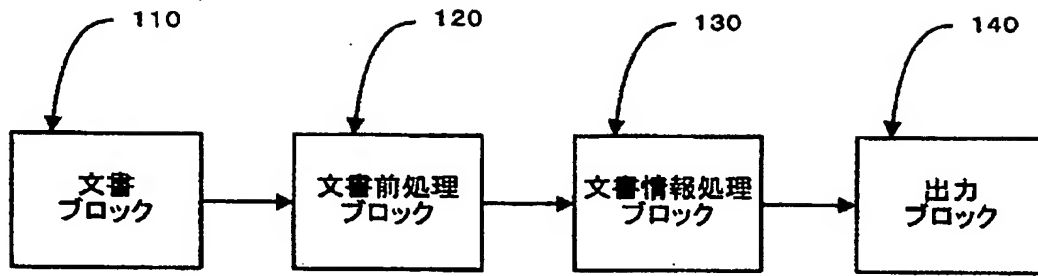
冗長なクラスターを削除する為のクラスターの重要度を説明する図である。

【符号の説明】

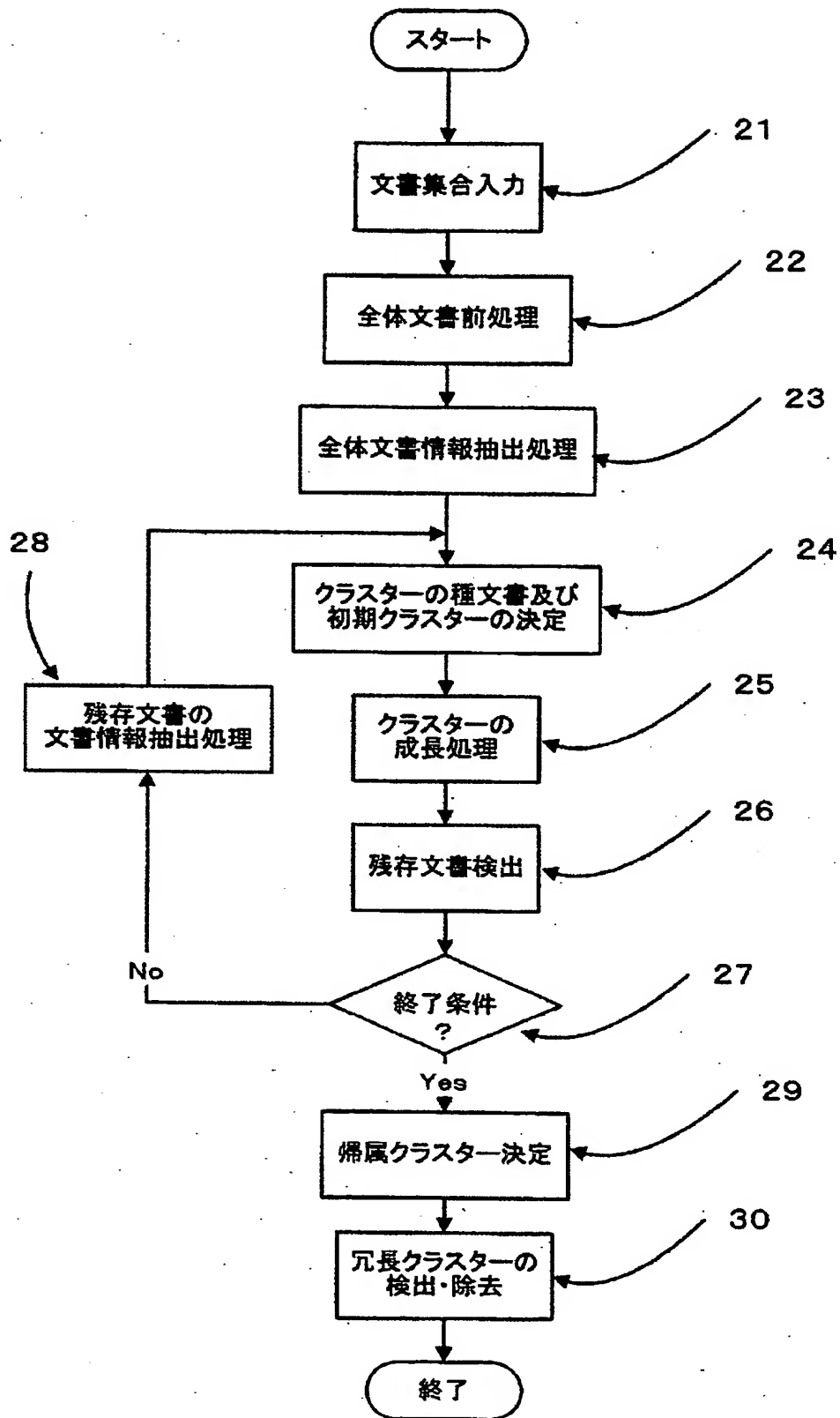
- 110 : 文書入力ブロック
- 120 : 文書前処理ブロック
- 130 : 文書情報処理ブロック
- 140 : 出力ブロック

【書類名】 図面

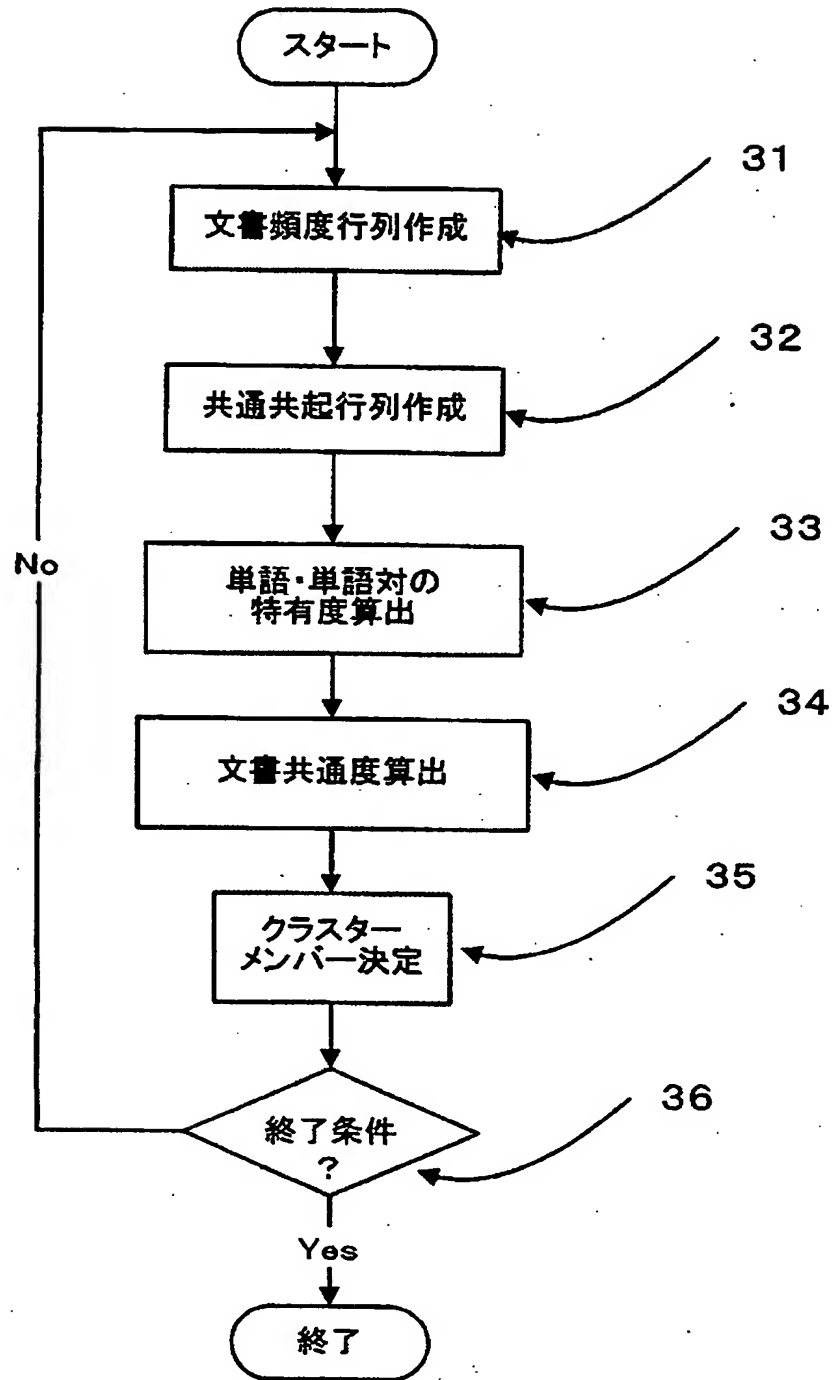
【図 1】



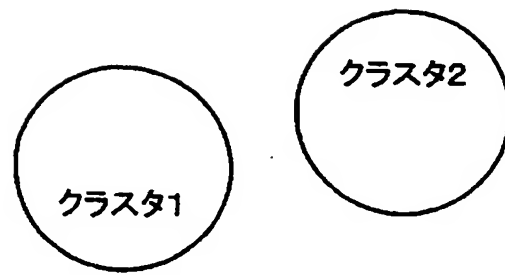
【図 2】



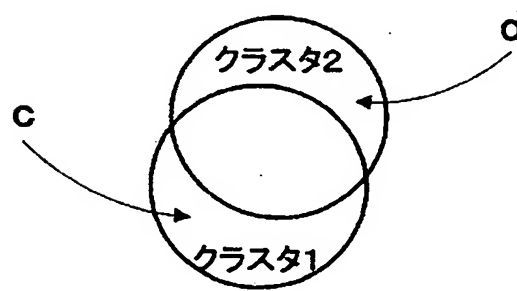
【図3】



【図 4】



(a)



(b)

【書類名】 要約書

【課題】

文書のクラスタリングにおいて、正しい数のクラスターを求めること、及び各文書の帰属するクラスターを精度よく求めることは完全には解決されていない問題であった。

【解決手段】

文書クラスタリングでは同じ話題を述べた文書がグループ化されるので、同じクラスターに属する文書群には何らかの共通性があるはずである。また、各話題には話題特有の用語や用語対が存在する。本発明ではこれらの点に着目し、各文書の着目クラスターへの近さを求めるときに、着目クラスターに特有でない用語や用語対の影響を排除しつつ着目クラスターの共通情報を用いるようにした。

【選択図】 図 1

特願 2003-105867

出 願 人 履 歴 情 報

識別番号

[398038580]

1. 変更年月日

1998年 5月19日

[変更理由]

新規登録

住 所

アメリカ合衆国カリフォルニア州パロアルト ハノーバー・ス
トリート 3000

氏 名

ヒューレット・パカード・カンパニー